



Generalization performance of regularized neural network models

Larsen, Jan; Hansen, Lars Kai

Published in:

Proceedings of the 4th IEEE Workshop Neural Networks for Signal Processing

Link to article, DOI:

[10.1109/NNSP.1994.366065](https://doi.org/10.1109/NNSP.1994.366065)

Publication date:

1994

Document Version

Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):

Larsen, J., & Hansen, L. K. (1994). Generalization performance of regularized neural network models. In *Proceedings of the 4th IEEE Workshop Neural Networks for Signal Processing* (pp. 42-51). IEEE.
<https://doi.org/10.1109/NNSP.1994.366065>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

GENERALIZATION PERFORMANCE OF REGULARIZED NEURAL NETWORK MODELS

Jan Larsen and Lars Kai Hansen
The Computational Neural Network Center
Electronics Institute, Building 349
Technical University of Denmark
DK-2800 Lyngby, Denmark

Abstract. Architecture optimization is a fundamental problem of neural network modeling. The optimal architecture is defined as the one which minimizes the generalization error. This paper addresses estimation of the generalization performance of regularized, complete neural network models. Regularization normally improves the generalization performance by restricting the model complexity. A formula for the optimal weight decay regularizer is derived. A regularized model may be characterized by an effective number of weights (parameters); however, it is demonstrated that no simple definition is possible. A novel estimator of the average generalization error (called *FPER*) is suggested and compared to the Final Prediction Error (*FPE*) and Generalized Prediction Error (*GPE*) estimators. In addition, comparative numerical studies demonstrate the qualities of the suggested estimator.

INTRODUCTION

One of the fundamental problems involved in design of neural network models is architecture optimization aiming at high generalization performance. In this paper the generalization measure is defined as the *average generalization error*, i.e., the expected squared error averaged over all possible training sets of size N , with N being the number of training samples. The average generalization error, Γ , can be decomposed into three additive components [2], [8]: $\Gamma = \sigma_\epsilon^2 + MSME + WFP$, viz. the inherent noise variance, the mean square model error, and the weight fluctuation penalty¹. The inherent noise variance is caused by noise on the data which – per definition – cannot be modeled.

¹The *MSME* and the *WFP* are related to the *squared bias* and the *variance*, respectively. See [2] for a definition of bias and variance.

Presence of *MSME* reflects the lack of modeling capability of the neural network for modeling the current data, i.e., the network is an *incomplete model* of the data generating system. Finally, the *WFP* reflects the increase in average generalization error caused by fluctuations in the estimated weights, which stem from the fact that the weights are estimated from a given finite training set.

Architecture optimization can be viewed as a bias/variance trade off [2], [11] or equivalently a *MSME/WFP* trade off: The *MSME* is reduced when increasing the network complexity² while the *WFP* typically³ increases. The literature provides a variety of methods for performing this trade off, including architecture pruning and growing schemes, as well as regularization techniques.

TRAINING AND GENERALIZATION

Consider modeling the data generating system:

$$y(k) = g(\mathbf{x}(k)) + \varepsilon(k) \quad (1)$$

where k is the discrete time index, $y(k)$ is the scalar output signal, $g(\cdot)$ constitutes a nonlinear mapping of the p -dimensional input signal $\mathbf{x}(k)$ (column vector), and $\varepsilon(k)$ is an inherent noise signal.

Assumption 1 *The input signal $\mathbf{x}(k)$ is assumed to be a strongly mixing⁴ strictly stationary sequence and the inherent noise $\varepsilon(k)$ is assumed to be a strictly stationary sequence independent on the input, white, with zero mean, and finite variance, σ_ε^2 .*

The neural network model of the system in Eq. (1) is given by

$$y(k) = f(\mathbf{x}(k); \mathbf{w}) + e(k; \mathbf{w}) \quad (2)$$

where $f(\cdot; \mathbf{w})$ defines the mapping of the neural network parameterized by the m -dimensional weight vector \mathbf{w} , and $e(k; \mathbf{w})$ is the error signal.

Assumption 2 *The model is assumed **complete** [8, Def. 6.3], i.e., there exists a true weight vector, \mathbf{w}° , so as to*

$$\forall \mathbf{x} : f(\mathbf{x}; \mathbf{w}^\circ) \equiv g(\mathbf{x}) \quad (3)$$

In general, only little a priori knowledge of the data generating system is available, i.e., most neural network models are *incomplete*, which result in non-zero mean square model error. However, a multi-layer perceptron neural

²This statement is only true for nested families of network architectures. Moreover, *MSME* may remain unchanged when adding irrelevant complexity.

³It should be emphasized that it is possible to give simple examples where the *WFP* actually *decreases* when adding extra complexity [8, Ch. 6.3.4].

⁴Loosely speaking, i.e., the dependence of $\mathbf{x}(k)$ and $\mathbf{x}(k+\tau)$ vanishes as $|\tau| \rightarrow \infty$.

network with many hidden neurons is capable of approximating a large class of functions, thus *MSME* may be small relative to $\sigma_\epsilon^2 + WFP$, and the model may be regarded as *quasi-complete*. When dealing with cases where the complete model assumption is dubious, it is suggested to estimate the generalization performance by using the *GEN* estimator [7], [8].

Define the training set of N samples by $T = \{\mathbf{x}(k); \mathbf{y}(k)\}$, $k = 1, 2, \dots, N$. The model is estimated by minimizing a cost function being the sum of the usual mean square cost and a weight decay regularizer⁵:

$$C_N(\mathbf{w}) = S_N(\mathbf{w}) + \mathbf{w}^\top \mathbf{R} \mathbf{w} \quad (4)$$

where $S_N(\mathbf{w}) = N^{-1} \sum_{k=1}^N e^2(k; \mathbf{w}) = N^{-1} \sum_{k=1}^N [y(k) - f(\mathbf{x}(k); \mathbf{w})]^2$ is the mean square cost and \mathbf{R} is a $m \times m$ symmetric, positive semidefinite regularization matrix. Standard weight decay regularization is obtained by using $\mathbf{R} = \kappa \mathbf{I}$, where $\kappa \geq 0$ is the weight decay parameter and \mathbf{I} the identity matrix. The presented theory is not restricted to the chosen cost function, thus analogous results can be obtained when e.g., using log-likelihood cost functions and more general regularizers, $r(\mathbf{w}; \kappa)$, where $r(\cdot)$ is a regularization function parameterized by κ .

The weights of the estimated model are denoted the *estimated weights*, i.e.,

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} C_N(\mathbf{w}) \quad (5)$$

Also define the *expected cost function*:

$$C(\mathbf{w}) = E\{C_N(\mathbf{w})\} = E\{e^2(\mathbf{w})\} + \mathbf{w}^\top \mathbf{R} \mathbf{w} \quad (6)$$

where $E\{\cdot\}$ denotes expectation w.r.t. the joint input-output probability density function. Under mild regularity conditions (see e.g., [8, Ch. 5], [12]) $\lim_{N \rightarrow \infty} C_N(\mathbf{w}) = C(\mathbf{w})$, and the estimated weight vector $\hat{\mathbf{w}}$ becomes a consistent estimate ($N \rightarrow \infty$) of the *optimal weight vector*: $\mathbf{w}^* = \arg \min_{\mathbf{w}} C(\mathbf{w})$. Since the model is assumed complete \mathbf{w}^* is identical to \mathbf{w}° when omitting regularization. However, regularization imposes a bias of the optimal weights towards $\mathbf{0}$.

The *generalization error* of the estimated model is defined as the expected squared error on an test sample, $[\mathbf{x}; \mathbf{y}]$, independent on the training samples, i.e.,

$$G(\hat{\mathbf{w}}) = E\{e^2(\hat{\mathbf{w}})\} = E\{[y - f(\mathbf{x}; \hat{\mathbf{w}})]^2\} \quad (7)$$

It turns out (see e.g., the discussion in [8, Sec. 6.3.2]) that $G(\hat{\mathbf{w}})$ is not necessarily a reliable measure of the model quality since it depends on the *actual training set* through $\hat{\mathbf{w}}$. In addition, it is not possible to obtain estimates of $G(\hat{\mathbf{w}})$ without perfect knowledge of the joint input-output distribution. Hence, the appropriate model quality measure is the *average generalization error*, e.g., [8], [11]:

$$\Gamma = E_T\{G(\hat{\mathbf{w}})\} \quad (8)$$

⁵ \top denotes the transpose operator.

where $E_{\mathcal{T}}\{\cdot\}$ denotes expectation over all training sets with N samples. That is, averaging is w.r.t. fluctuation in $\hat{\mathbf{w}}$ due to different training sets. Define $\mathcal{T}_{\mathbf{x}} = \{\mathbf{x}(k)\}$ and $\mathcal{T}_{\epsilon} = \{\epsilon(k)\}$. As the noise and the input are assumed independent, the expectation w.r.t. \mathcal{T} is carried out as⁶:

$$E_{\mathcal{T}}\{G\} = E_{\mathcal{T}_{\mathbf{x}}}\{E_{\mathcal{T}_{\epsilon}}\{G|\mathcal{T}_{\mathbf{x}}\}\} \quad (9)$$

ESTIMATING THE AVERAGE GENERALIZATION ERROR

The objective of this presentation is to obtain an estimate of Γ defined in Eq. (8) calculated in terms of quantities derived from the estimated model. From a statistical point of view it is possible to set different quality requirements on the estimator. Here the following requirements are made:

Definition 1 *The estimator searched for, $\hat{\Gamma}$, is required to be consistent, and unbiased to order $1/N$, i.e., $\hat{\Gamma} \rightarrow \Gamma$ as $N \rightarrow \infty$, and $E_{\mathcal{T}}\{\hat{\Gamma}\} = \Gamma + o(1/N)$, where $o(\cdot)$ is the order function.*

The basic tool for deriving an estimator are second order Taylor series expansions of the average training and generalization errors, as follows:

$$\begin{aligned} E_{\mathcal{T}}\{S_N(\hat{\mathbf{w}})\} &\approx E_{\mathcal{T}}\{S_N(\mathbf{w}^{\circ})\} + E_{\mathcal{T}}\left\{\frac{\partial S_N(\mathbf{w}^{\circ})}{\partial \mathbf{w}^{\top}} \Delta \mathbf{w}\right\} \\ &\quad + E_{\mathcal{T}}\{\Delta \mathbf{w}^{\top} \mathbf{H}_N(\mathbf{w}^{\circ}) \Delta \mathbf{w}\} \end{aligned} \quad (10)$$

$$\begin{aligned} E_{\mathcal{T}}\{G(\hat{\mathbf{w}})\} &\approx E_{\mathcal{T}}\{G(\mathbf{w}^{\circ})\} + E_{\mathcal{T}}\left\{\frac{\partial G(\mathbf{w}^{\circ})}{\partial \mathbf{w}^{\top}} \Delta \mathbf{w}\right\} \\ &\quad + E_{\mathcal{T}}\{\Delta \mathbf{w}^{\top} \mathbf{H}(\mathbf{w}^{\circ}) \Delta \mathbf{w}\} \end{aligned} \quad (11)$$

where $\Delta \mathbf{w}$ is the weight fluctuation $\Delta \mathbf{w} = \hat{\mathbf{w}} - \mathbf{w}^{\circ}$, $\mathbf{H}_N(\mathbf{w})$ is the Hessian matrix of the mean square cost function, i.e.,

$$\mathbf{H}_N(\mathbf{w}) = \frac{1}{2} \frac{\partial^2 S_N(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^{\top}} = \frac{1}{N} \sum_{k=1}^N \psi(k; \mathbf{w}) \psi^{\top}(k; \mathbf{w}) - \Psi(k; \mathbf{w}) e(k; \mathbf{w}) \quad (12)$$

defining ψ as the instantaneous gradient vector of the model output, $\psi(k; \mathbf{w}) = \partial f(\mathbf{x}(k); \mathbf{w}) / \partial \mathbf{w}$. Finally, Ψ is the second derivative matrix of the model output, $\Psi(k; \mathbf{w}) = \partial \psi(k; \mathbf{w}) / \partial \mathbf{w}^{\top}$. Similarly, $\mathbf{H}(\mathbf{w})$ is the Hessian matrix of the generalization error, given by

$$\mathbf{H}(\mathbf{w}) = \frac{1}{2} \frac{\partial^2 G(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^{\top}} = E\left\{\psi(\mathbf{w}) \psi^{\top}(\mathbf{w}) - \Psi(\mathbf{w}) e(\mathbf{w})\right\} \quad (13)$$

In order to ensure the validity of the Taylor series approximations it is required that $\Delta \mathbf{w}$ is sufficiently small. As mentioned above $\hat{\mathbf{w}}$ is a consistent estimate

⁶Note that expectation over the training set, $\mathcal{T} = \{\mathbf{x}(k); y(k)\}$, equals expectation over input and inherent noise samples, cf. the model definition Eq. (2).

of \mathbf{w}^* ; however, \mathbf{w}^* does not collapse onto \mathbf{w}^0 unless $\mathbf{R} = \mathbf{0}$. Consequently, it is expected that the Taylor series are valid for sufficiently large N and sufficiently small \mathbf{R} .

The appendix below provides a brief evaluation of the individual terms of Eq. (10), (11). The result is: For $N > 2m_1 - m_2$,

$$E_{\mathcal{T}} \{S_N(\hat{\mathbf{w}})\} = \sigma_{\varepsilon}^2 \left(1 - \frac{2m_1 - m_2}{N}\right) + M' + o(1/N) \quad (14)$$

$$\Gamma = \sigma_{\varepsilon}^2 \left(1 + \frac{m_2}{N}\right) + M + o(1/N) \quad (15)$$

where m_1, m_2 defines two different *effective number of weights*⁷:

$$m_1 = \text{tr} [\mathbf{H}(\mathbf{w}^0) \mathbf{J}^{-1}(\mathbf{w}^0)], \quad m_2 = \text{tr} [\mathbf{H}(\mathbf{w}^0) \mathbf{J}^{-1}(\mathbf{w}^0) \mathbf{H}(\mathbf{w}^0) \mathbf{J}^{-1}(\mathbf{w}^0)] \quad (16)$$

$\mathbf{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \mathbf{R}$ is the Hessian matrix of the expected cost function which is assumed to be invertible, and $\text{tr}[\cdot]$ is the trace operator.

$$M' = (\mathbf{w}^0)^{\top} \mathbf{R} \mathbf{J}^{-1}(\mathbf{w}^0) \left(\mathbf{H}(\mathbf{w}^0) + \frac{-2\mathbf{K}_1 + \mathbf{K}_2}{N} \right) \mathbf{J}^{-1}(\mathbf{w}^0) \mathbf{R} \mathbf{w}^0 \quad (17)$$

with $\mathbf{K}_1, \mathbf{K}_2$ being 4th order moments, as shown by⁸:

$$\mathbf{K}_1 = E \left\{ \left(\boldsymbol{\psi} \boldsymbol{\psi}^{\top} - \mathbf{H} \right) \mathbf{J}^{-1} \left(\boldsymbol{\psi} \boldsymbol{\psi}^{\top} - \mathbf{H} \right) \right\} \quad (18)$$

$$\mathbf{K}_2 = E \left\{ \left(\boldsymbol{\psi} \boldsymbol{\psi}^{\top} - \mathbf{H} \right) \mathbf{J}^{-1} \mathbf{H} \mathbf{J}^{-1} \left(\boldsymbol{\psi} \boldsymbol{\psi}^{\top} - \mathbf{H} \right) \right\} \quad (19)$$

M equals M' except that the term \mathbf{K}_1 is absent. In general, M and M' are negligible compared to the remaining terms in Eq. (14), (15) when 1) using a regularization matrix close to the optimal setting Eq. (24), and when 2) the signal-to-noise ratio, $V\{g(\mathbf{x})\}/\sigma_{\varepsilon}^2$, is reasonable large.

Neglecting M, M' and eliminating σ_{ε}^2 in Eq. (14), (15) leads to:

$$\hat{\Gamma} = \frac{N + m_2}{N - 2m_1 + m_2} E_{\mathcal{T}} \{S_N(\hat{\mathbf{w}})\}, \quad N > 2m_1 - m_2 \quad (20)$$

which is unbiased to $o(1/N)$. Notice that elimination of σ_{ε}^2 introduces terms proportional to N^{-j} , $j > 1$. This seems inconsistent; however, for practical purposes the form is convenient since $\hat{\Gamma}$ typically is an underestimate of Γ on the average. In the case of a complete linear model which is estimated without regularization [3] and [8, Theorem 6.10] support this statement.

The suggested estimator may be viewed as an extension of the classical *FPE* estimator [1], $FPE = E_{\mathcal{T}} \{S_N(\hat{\mathbf{w}})\} (N + m)/(N - m)$, in which the

⁷It is easily shown that $m_1 \geq m_2 > 0$ thus $2m_1 - m_2 > 0$. Moreover, 1) $m_1 = m_2 = m$ for $\mathbf{R} = \mathbf{0}$ and $\mathbf{H}(\mathbf{w}^0)$ non-singular, and 2) $m_1 \rightarrow 0, m_2 \rightarrow 0$ as $\|\mathbf{R}\| \rightarrow \infty$.

⁸ $\mathbf{K}_1, \mathbf{K}_2$ are of order one, and limited by assumption. Further note that all involved quantities are evaluated at \mathbf{w}^0 .

number of weights m is replaced by the different effective number of weights, m_2 and $2m_1 - m_2$. Moreover, the estimator can be interpreted as a special version⁹ of the *GPE* estimator [10], [11] where the inherent noise variance is estimated by: $\sigma_\epsilon^2 = E_{\mathcal{T}} \{S_N(\hat{\mathbf{w}})\} N/(N - 2m_1 + m_2)$. In order to construct a Γ -estimator from observable quantities, estimation of the noise variance is indeed important. This problem is not directly addressed in [10], [11]. The estimator suggested in [10] reads: $\sigma_\epsilon^2 = E_{\mathcal{T}} \{S_N(\hat{\mathbf{w}})\} N/(N - m_1)$, which obviously differs from the one derived from Eq. (14). In conclusion – as suggested in [9], [11] – it is not possible to define a single quantity m_1 which expresses the effective number of weights in the model, since σ_ϵ^2 should be estimated from $2m_1 - m_2$ rather than m_1 effective weights.

For practical purposes the quantities in Eq. (20) are estimated from observed quantities. An unbiased $o(1/N)$ estimator within the second order Taylor series expansion Eq. (10), (11) is the the **Final Prediction Error** estimator for **R**egularized models, as shown by:

$$FPER = \frac{N + \hat{m}_2}{N - 2\hat{m}_1 + \hat{m}_2} S_N(\hat{\mathbf{w}}), \quad N > 2\hat{m}_1 - \hat{m}_2 \quad (21)$$

where

$$\hat{m}_1 = \text{tr} [\mathbf{H}_N(\hat{\mathbf{w}}) \mathbf{J}_N^{-1}(\hat{\mathbf{w}})], \quad \hat{m}_2 = \text{tr} [\mathbf{H}_N(\hat{\mathbf{w}}) \mathbf{J}_N^{-1}(\hat{\mathbf{w}}) \mathbf{H}_N(\hat{\mathbf{w}}) \mathbf{J}_N^{-1}(\hat{\mathbf{w}})] \quad (22)$$

and $\mathbf{J}_N(\hat{\mathbf{w}}) = \mathbf{H}_N(\hat{\mathbf{w}}) + \mathbf{R}$ is the Hessian matrix of the cost function which is assumed to be invertible.

OPTIMIZING THE WEIGHT DECAY REGULARIZATION PARAMETER

For simplicity, consider simple weight decay regularization, i.e., $\mathbf{R} = \kappa \mathbf{I}$ where κ is the weight decay parameter. As mentioned in the introduction, trading off weight fluctuation penalty (*WFP*) and mean square model error (*MSME*) leads to an optimal setting of κ . In [6] this problem was addressed for linear models and the following may be viewed as an extension of this work.

Inspecting Eq. (15) it turns out that¹⁰ $M = MSME$ and $WFP = \sigma_\epsilon^2 m_2/N$. The optimal value, κ_{opt} , is found by solving:

$$\frac{\partial WFP}{\partial \kappa} + \frac{\partial MSME}{\partial \kappa} = 0 \quad (23)$$

As expected, $\lim_{N \rightarrow \infty} WFP = 0$, since it measures the contribution due to a finite training set. Consequently, in order to reach the minimal average generalization error $\Gamma = \sigma_\epsilon^2$ the restriction $\lim_{N \rightarrow \infty} MSME = 0$ should be met. The κ -dependence of the individual elements of \mathbf{K}_1 is $(\lambda_i + \kappa)^{-1}$

⁹Notice that this coincidence is based on various important assumptions, e.g., the model being complete and the negligibility of M .

¹⁰Notice when determining an optimal κ , M is *not* neglected in Eq. (15).

where λ_i is the i 'th eigenvalue of $\mathbf{H}(\mathbf{w}^\circ)$. For \mathbf{K}_2 the element dependence is: $\prod_{i=1}^2 (\lambda_{i_r} + \kappa)^{-1}$. In summary, M is a sum of addends which κ -dependence are given by: $\kappa^2 \prod_{r=1}^\ell (\lambda_{i_r} + \kappa)^{-1}$, $\ell \in \{2, 3, 4\}$. That is, to fulfill the requirement $\lim_{N \rightarrow \infty} MSME = 0$, $\lim_{N \rightarrow \infty} \kappa = 0$ should be imposed. The solution to Eq. (23) can therefore be expressed as: $\kappa_{\text{opt}} = \kappa'_{\text{opt}}/N + o(1/N)$. Expanding the addends of Eq. (23) to first order in κ and $1/N$ and solving for κ gives:

$$\kappa_{\text{opt}} = \frac{\sigma_\epsilon^2}{N} \cdot \frac{\text{tr } \mathbf{H}^+(\mathbf{w}^\circ)}{(\mathbf{w}^\circ)^\top \mathbf{H}^+(\mathbf{w}^\circ) \mathbf{w}^\circ} + o(1/N) \quad (24)$$

where $\mathbf{H}^+(\mathbf{w}^\circ)$ is the Moore-Penrose pseudo inverse. Suppose the eigenvalues of $\mathbf{H}(\mathbf{w}^\circ)$ obey: $\lambda_1 \geq \dots \geq \lambda_n > 0$ and $\lambda_i = 0$, $\forall i \in [n+1; m]$. The associated eigenvectors are assembled (as column vectors) in the matrix \mathbf{Q} . The pseudo inverse then reads: $\mathbf{H}^+(\mathbf{w}^\circ) = \mathbf{Q} \text{diag}[\lambda_1^{-1}, \dots, \lambda_n^{-1}, 0, \dots, 0] \mathbf{Q}^\top$.

Notice two facts concerning κ_{opt} : First, it is proportional to the inherent noise variance. If no noise is present $WFP = 0$, thus one should not introduce $MSME$ by employing a non-zero κ . Secondly, κ_{opt} is inversely proportional to the length of the optimal weight vector weighted by the elements of the Moore-Penrose pseudo inverse Hessian matrix. This is due to the fact that we regularize against the zero weight vector.

Since the optimal weights \mathbf{w}° are unknown, it is impossible to calculate κ_{opt} directly; however, in [4] *adaptive regularization* is studied for a linear one-dimensional model, and [5] presents an adaptive regularization scheme for the purpose of designing compact time series models. In addition, it is possible to show that the average generalization error is reduced when using $0 < \kappa \leq 2\kappa_{\text{opt}}$.

NUMERICAL EXPERIMENTS

To substantiate the qualities of the suggested *FPER* estimator Eq. (21), numerical comparisons with the *FPE* and *GPE* estimators¹¹,

$$FPE = \frac{N+m}{N-m} S_N(\hat{\mathbf{w}}) \quad GPE = \frac{N+\hat{m}_1}{N-\hat{m}_1} S_N(\hat{\mathbf{w}}) \quad (25)$$

is – for convenience – performed for a linear model. The linear data generating system (dimension $m = 15$) is given by:

$$\mathbf{y}(k) = \mathbf{x}^\top(k) \mathbf{w}^\circ + \epsilon(k) \quad (26)$$

where $\mathbf{x}(k)$ is an i.i.d. Gaussian distributed sequence with zero mean and, the elements of $\mathbf{H} = E\{\mathbf{x}\mathbf{x}^\top\}$ are selected randomly, resulting in an eigenvalue-spread approx. equal to 900. The optimal weights are drawn independently from a standard Gaussian distribution. The inherent noise is a Gaussian

¹¹As regards the *GPE* estimator, the noise variance estimation suggested in [10] is employed.

zero mean, i.i.d. sequence which is independent of the input with variance $\sigma_\varepsilon^2 = 0.25 \cdot E \left\{ (\mathbf{x}^\top(\mathbf{k})\mathbf{w}^\circ)^2 \right\} = 0.25 \cdot (\mathbf{w}^\circ)^\top \mathbf{H} \mathbf{w}^\circ$. That is, the signal-to-noise ratio equals approx. 6 dB.

$Q = 2.4 \cdot 10^4$ independent training sets of size N in the interval $[15; 35]$ were randomly generated, and the weights of the associated model were estimated using a simple weight decay regularizer with $\kappa = 2\kappa_{\text{opt}}$.

The “true” average generalization error was estimated by $\hat{\Gamma}_G = \langle G(\hat{\mathbf{w}}) \rangle$ where $\langle \cdot \rangle$ denotes the average w.r.t. the Q training sets, and

$$G(\hat{\mathbf{w}}) = E \left\{ [\varepsilon + \mathbf{x}^\top (\mathbf{w}^\circ - \hat{\mathbf{w}})]^2 \right\} = \sigma_\varepsilon^2 + (\mathbf{w}^\circ - \hat{\mathbf{w}})^\top \mathbf{H} (\mathbf{w}^\circ - \hat{\mathbf{w}}) \quad (27)$$

The quality of the estimators¹², $\hat{\Gamma}(T) \in \{FPER, FPE, GPE\}$, is quantified by three different measures:

$$NB = \frac{\hat{\Gamma}(T) - \hat{\Gamma}_G}{\hat{\Gamma}_G} \cdot 100\% \quad NRMSE = \frac{\sqrt{\left\langle [\hat{\Gamma}(T) - \hat{\Gamma}_G]^2 \right\rangle}}{\hat{\Gamma}_G} \cdot 100\% \quad (28)$$

$$\Pi(\hat{\Gamma}) = \left\langle \mu \left(\left| \hat{\Gamma}(T) - \hat{\Gamma}_G \right| - \left| FPER(T) - \hat{\Gamma}_G \right| \right) \right\rangle \quad (29)$$

NB is the normalized bias, $NRMSE$ is the normalized root mean square error, and Π is the probability that $FPER$ is closer to the true estimate, $\hat{\Gamma}_G$, than another estimator, $\hat{\Gamma}$. Here $\mu(\cdot)$ denotes the step function. Fig. 1 shows plots of the considered measures. NB of $FPER$ is smallest for all training set sizes; however, as the training set size approaches infinity all estimates becomes identical as $\kappa_{\text{opt}} \rightarrow 0$. For $N = 35$ $NB(FPER)$ is approx. half the $NB(GPE)$. The $NRMSE$'s of $FPER$ and GPE are approx. identical, thus one could claim that the normalized bias improvement of $FPER$ relative to GPE is lost at increased variance¹³. However, the probability that $FPER$ is closer than GPE to the true Γ is around 0.65; consequently, $FPER$ should be preferred to GPE . FPE shows extremely bad performance in all figures and moreover, FPE is *negative*, possibly infinite when $N \leq 15$.

CONCLUSION

This paper presented an consistent and $o(1/N)$ unbiased estimator of the average generalization error for a complete neural network model, called $FPER$. The network is trained by using a cost function which is the sum of the mean square error and a quadratic regularization term. The estimator may be viewed as an extension of the FPE and GPE estimators [1], [10]. It turns out that the complexity reduction obtained by using regularization is expressed in terms of *two* distinct effective number of weights, unlike defining a single

¹²Notice, the dependence on the particular training set, T , is emphasized.

¹³That is, mean square error the minus squared bias.

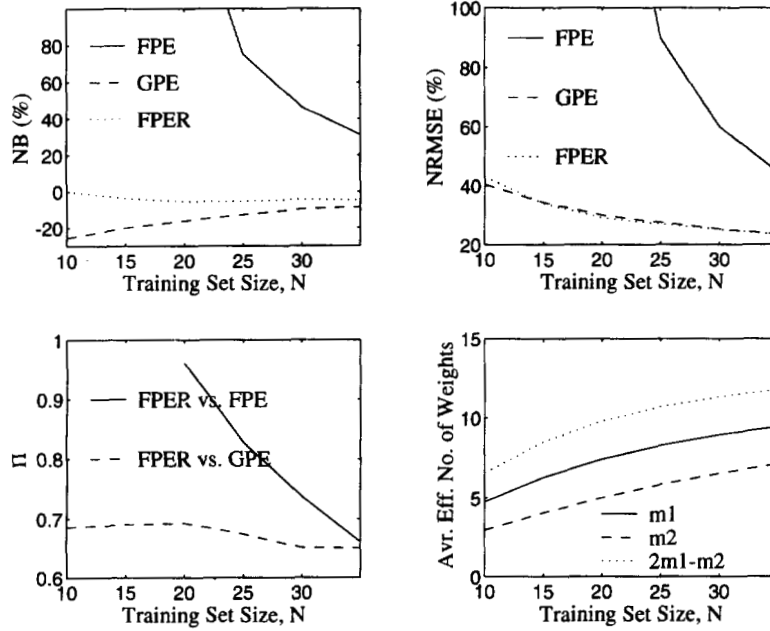


Figure 1: Comparison of *FPE*, *GPE* and *FPER*. The *FPE* curves are not calculated for $N \leq 15$, and the upper panels are cutoff at +100%. The bottom right panel shows the average effective number of weights $\langle \hat{m}_1 \rangle$, $\langle \hat{m}_2 \rangle$ as well as $2\langle \hat{m}_1 \rangle - \langle \hat{m}_2 \rangle$, quantity reflecting the effective number of weights, as suggested in [9], [11]. Moreover, an expression for the optimal weight decay parameter is presented and discussed. The potential of the *FPER* estimator was demonstrated by comparative numerical studies.

ACKNOWLEDGMENTS

This work was supported by the Danish Natural Science and Technical Research Councils through the Computational Neural Network Center.

APPENDIX

Evaluation of the terms in Eq. (10), (11) is based on two observations: First, $\partial C_N(\hat{\mathbf{w}})/\partial \mathbf{w} = 0$ since $\hat{\mathbf{w}}$ minimizes $C_N(\mathbf{w})$. A first order Taylor series expansion of $\partial C_N(\hat{\mathbf{w}})/\partial \mathbf{w}$ reads¹⁴:

$$\frac{\partial C_N(\mathbf{w}^0)}{\partial \mathbf{w}} + \frac{\partial^2 C_N(\mathbf{w}^0)}{\partial \mathbf{w} \partial \mathbf{w}^\top} \Delta \mathbf{w} = 0 \quad (30)$$

¹⁴Expanding beyond first order result in 3rd and higher order derivatives of the cost function which already are assumed to be negligible.

Subsequently, a few algebraic manipulations result in:

$$\Delta \mathbf{w} = \mathbf{J}_N^{-1}(\mathbf{w}^o) \left[\frac{1}{N} \sum_{k=1}^N \psi(k; \mathbf{w}^o) \varepsilon(k) - \mathbf{R} \mathbf{w}^o \right] \quad (31)$$

where $\mathbf{J}_N(\mathbf{w}^o)$ is the non-singular Hessian matrix of the cost function.

The second observation is an expansion of the inverse Hessian obtained by repeatedly using the matrix inversion lemma [8, App. A,B]. The result is: $\mathbf{J}_N^{-1}(\mathbf{w}^o) = \mathbf{J}^{-1}(\mathbf{w}^o) - N^{-1} \cdot \mathbf{J}^{-1}(\mathbf{w}^o) \boldsymbol{\Theta} \mathbf{J}^{-1}(\mathbf{w}^o) + \dots$, where $\boldsymbol{\Theta} = \mathbf{H}_N(\mathbf{w}^o) - \mathbf{H}(\mathbf{w}^o)$.

REFERENCES

- [1] H. Akaike, "Fitting Autoregressive Models for Prediction," Annals of the Institute of Statistical Mathematics, vol. 21, pp. 243-247, 1969.
- [2] S. Geman, E. Bienenstock & R. Doursat, "Neural Networks and the Bias/Variance Dilemma," Neural Computation, vol. 4, pp. 1-58, 1992.
- [3] L.K. Hansen, "Stochastic Linear Learning: Exact Test and Training Error Averages," Neural Networks, vol. 6, pp. 393-396, 1993.
- [4] L.K. Hansen & C.E. Rasmussen, "Pruning from Adaptive Regularization," Preprint Electronics Institute, The Technical University of Denmark, 1993. Accepted for publication in Neural Computation.
- [5] L.K. Hansen, C.E. Rasmussen, C. Svarer, & J. Larsen, "Adaptive Regularization," in Proceedings of the 1994 IEEE NNSP Workshop.
- [6] A. Krogh & J.A. Hertz, "A Simple Weight Decay Can Improve Generalization," in J.E. Moody, S.J. Hanson, R.P. Lippmann (eds.) Advances in Neural Information Processing Systems 4, Proceedings of the 1991 Conference, San Mateo, California: Morgan Kaufmann Publishers, 1992, pp. 950-957.
- [7] J. Larsen, "A Generalization Error Estimate for Nonlinear Systems," in S.Y. Kung, F. Fallside, J. Aa. Sørensen & C.A. Kamm (eds.) Neural Networks for Signal Processing 2: Proceedings of the 1992 IEEE-SP Workshop, Piscataway, New Jersey: IEEE, 1992, pp. 29-38.
- [8] J. Larsen, Design of Neural Network Filters, Ph.D. Thesis, Electronics Institute, The Technical University of Denmark, March 1993.
- [9] D. MacKay, "A Practical Bayesian Framework for Backprop Networks," Neural Computation, vol. 4, pp. 448-472, 1992.
- [10] J. Moody, "Note on Generalization, Regularization, and Architecture Selection in Nonlinear Learning Systems," in B.H. Juang, S.Y. Kung & C.A. Kamm (eds.) Proceedings of the first IEEE Workshop on Neural Networks for Signal Processing, Piscataway, New Jersey: IEEE, 1991, pp. 1-10.
- [11] J. Moody, "The Effective Number of Parameters: An Analysis of Generalization and Regularization in Nonlinear Learning Systems," in J.E. Moody, S.J. Hanson, R.P. Lippmann (eds.) Advances in Neural Information Processing Systems 4, Proceedings of the 1991 Conference, San Mateo, California: Morgan Kaufmann Publishers, 1992, pp. 847-854.
- [12] H. White, "Consequences and Detection of Misspecified Nonlinear Regression Models," Journal of the American Statistical Association, vol. 76, no. 374, pp. 419-433, June 1981.

AN APPLICATION OF IMPORTANCE-BASED FEATURE EXTRACTION IN REINFORCEMENT LEARNING

David J. Finton
Computer Sciences Department
University of Wisconsin-Madison
Madison, WI 53706

Yu Hen Hu
Department of Electrical and Computer Engineering
University of Wisconsin-Madison
Madison, WI 53706

Abstract—The sparse feedback in reinforcement learning problems makes feature extraction difficult. We present *importance-based feature extraction*, which guides a bottom-up self-organization of feature detectors according to top-down information as to the importance of the features; we define importance in terms of the reinforcement values expected as a result of taking different actions when a feature is recognized. We illustrate these ideas in terms of the pole-balancing task and a learning system which combines bottom-up tuning with a distributed version of Q-learning; adding importance-based feature extraction to the detector tuning resulted in faster learning.

INTRODUCTION

In reinforcement learning problems the feedback is simply a scalar value which may be delayed in time. This reinforcement signal reflects the success or failure of the entire system after it has performed some sequence of actions. Hence the reinforcement signal does not assign credit or blame to any one action (the *temporal credit assignment* problem), or to any particular node or system element (the *structural credit assignment* problem).

Since the reinforcement feedback is not an error signal for individual system elements, it gives little guidance for feature extraction, the on-line development of the system's input representation. Acting properly depends on both identifying the current context as well as selecting an action appropriate to that context, but the scalar feedback signal does not indicate which of these processes is at fault. It does not indicate whether the system should tune its feature detectors, or the weights placed on the outputs of those feature detectors, or both.